

Matthew Spahl, Emerson Kiefer, Kevin Martell Luya, Owen Tibby

CS 532 Final Project Instructions for running code

(also providing as a txt file titled “README”, but this is a bit longer and allows for pictures)

This program reads in a dataset of used car pricing information, trains a random forest on that dataset, and then predicts the RV_percent (retained value percent) for each sample in the test data.

Code is contained in the file Pyspark_forest.py

Download the dataset below and put it in the same directory as Pyspark_forest.py. The data file should be named: cleaned_data_CS532.csv

Link to data set that has been further cleaned (the one to use with the program):

https://drive.google.com/file/d/1X4a0lvQ9_nXYGUkL6uiyy5Zvu2wYlhEP/view?usp=sharing

Original Dataset from Kaggle (for reference but not the one to use with the program):

<https://www.kaggle.com/datasets/tunguz/used-car-auction-prices/data>

References consulted to learn how to use a random forest classifier in PySpark and measure time elapsed:

<https://towardsdatascience.com/a-guide-to-exploit-random-forest-classifier-in-pyspark-46d6999cb5db>

<https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.ml.regression.RandomForestRegressor.html>

<https://www.machinelearningplus.com/pyspark/pyspark-onehot-encoding/>

[https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.ml.evaluation.RegressionEvaluator.html](https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.ml.evaluation.ReggressionEvaluator.html)

<https://stackoverflow.com/questions/7370801/how-do-i-measure-elapsed-time-in-python>

Instructions for running code: in an Ubuntu VM with python and pip installed, install numpy and pandas.

To run, inside the directory with the Pyspark_forest.py file, use command: python3 Pyspark_forest.py

By default, the program will train/fit the random forest model using 50 trees. To change this, either change the value of treeCount on line 110 of the program or provide a command line argument. For example, to run with 100 trees, use the command:

python3 Pyspark_forest.py 100

The program will process the data and train a random forest regressor model and make predictions, then output the total runtime of the program, the data processing time, and the time for train/fit the random forest, make predictions, and evaluate the model.

Also, in our analysis for our presentation, referenced spark.apache.org/docs/2.2.0/mllib-ensembles.html to confirm that Pyspark Random Forest is able to train trees in parallel.

Example output from 6GB RAM and 6 cores, 200 trees:

```

matt@ubuntu-vm-2: ~/Documents/cs532/Final_Project
23/12/01 15:38:08 WARN MemoryStore: Not enough space to cache rdd_80_0 in memory! (computed 19.9 MiB so far)
23/12/01 15:38:08 WARN BlockManager: Persisting block rdd_80_0 to disk instead.
23/12/01 15:38:22 WARN DAGScheduler: Broadcasting large task binary with size 1345.0 KiB
23/12/01 15:38:28 WARN DAGScheduler: Broadcasting large task binary with size 2.5 MiB
23/12/01 15:38:37 WARN DAGScheduler: Broadcasting large task binary with size 4.9 MiB
23/12/01 15:38:51 WARN DAGScheduler: Broadcasting large task binary with size 1222.0 KiB
+-----+
|RV_percent| prediction|
+-----+
16.8| 17.47097039416918|
24.2| 21.19060194137648|
24.9| 22.175185499558324|
27.1| 28.31705232171758|
22.6| 20.69620922874034|
29.0| 28.307127592492648|
21.3| 20.77592610508753|
19.9| 20.77592610508753|
32.7| 28.638854012353715|
24.6| 23.612870165846758|
34.4| 35.446971016194716|
20.2| 16.04234322207186|
22.1| 21.228771698422307|
30.1| 21.228771698422307|
18.7| 16.04234322207186|
19.4| 17.334976842491386|
25.9| 22.433238244055246|
19.4| 19.1157642706606|
15.9| 19.1157642706606|
28.4| 24.42066372031255|
20.6| 21.158760967355747|
23.9| 22.424428867078223|
22.0| 24.59301633048278|
16.3| 20.243576712740136|
16.1| 18.46302905900124|
20.2| 22.55081343165531|
26.9| 27.485180384671036|
16.1| 20.187159864814873|
25.2| 23.539247725083285|
16.0| 22.408594043245756|
+-----+
only showing top 30 rows

RMSE: 9.441026787012408
Total time elapsed: 90.2157244682312
Time to preprocess: 35.730706214904785
Time to train, predict, and evaluate: 54.48501753807068
matt@ubuntu-vm-2:~/Documents/cs532/Final_Project$

```

As explained in our presentation, when using 2 or 4 cores, works for up to 400 trees, but with 6 cores, worked for up to 200 trees, then get the following out of memory error if we try to use 300 trees (we document this in our presentation, this is due to PySpark using parallelism with more cores and trying to train multiple trees at once):

```

matt@ubuntu-vm-2: ~/Documents/cs532/Final_Project
23/12/01 15:48:59 WARN MemoryStore: Not enough space to cache rdd_80_4 in memory! (computed 68.1 MiB so far)
23/12/01 15:49:04 WARN MemoryStore: Not enough space to cache rdd_80_0 in memory! (computed 30.0 MiB so far)
23/12/01 15:49:04 WARN MemoryStore: Not enough space to cache rdd_80_4 in memory! (computed 45.2 MiB so far)
23/12/01 15:49:04 WARN MemoryStore: Not enough space to cache rdd_80_3 in memory! (computed 68.1 MiB so far)
23/12/01 15:49:10 WARN DAGScheduler: Broadcasting large task binary with size 1012.6 KiB
23/12/01 15:49:10 WARN MemoryStore: Not enough space to cache rdd_80_4 in memory! (computed 30.0 MiB so far)
23/12/01 15:49:10 WARN MemoryStore: Not enough space to cache rdd_80_0 in memory! (computed 68.1 MiB so far)
23/12/01 15:49:10 WARN MemoryStore: Not enough space to cache rdd_80_3 in memory! (computed 45.2 MiB so far)
23/12/01 15:49:19 WARN DAGScheduler: Broadcasting large task binary with size 1932.0 KiB
23/12/01 15:49:20 WARN MemoryStore: Not enough space to cache rdd_80_4 in memory! (computed 30.0 MiB so far)
23/12/01 15:49:20 WARN MemoryStore: Not enough space to cache rdd_80_3 in memory! (computed 68.1 MiB so far)
23/12/01 15:49:20 WARN MemoryStore: Not enough space to cache rdd_80_0 in memory! (computed 45.2 MiB so far)
23/12/01 15:49:33 WARN DAGScheduler: Broadcasting large task binary with size 3.7 MiB
23/12/01 15:49:34 WARN MemoryStore: Not enough space to cache rdd_80_4 in memory! (computed 45.2 MiB so far)
23/12/01 15:49:34 WARN MemoryStore: Not enough space to cache rdd_80_3 in memory! (computed 30.0 MiB so far)
23/12/01 15:49:34 WARN MemoryStore: Not enough space to cache rdd_80_0 in memory! (computed 68.1 MiB so far)
23/12/01 15:49:58 WARN DAGScheduler: Broadcasting large task binary with size 7.2 MiB
23/12/01 15:49:58 WARN MemoryStore: Not enough space to cache rdd_80_3 in memory! (computed 45.2 MiB so far)
23/12/01 15:49:58 WARN MemoryStore: Not enough space to cache rdd_80_5 in memory! (computed 45.2 MiB so far)
23/12/01 15:49:58 WARN MemoryStore: Not enough space to cache rdd_80_0 in memory! (computed 45.2 MiB so far)
23/12/01 15:49:59 WARN MemoryStore: Not enough space to cache rdd_80_4 in memory! (computed 45.2 MiB so far)
23/12/01 15:50:09 ERROR Executor: Exception in task 3.0 in stage 36.0 (TID 144)
java.lang.OutOfMemoryError: GC overhead limit exceeded
    at java.lang.Double.valueOf(Double.java:519)
    at scala.runtime.BoxesRuntime.boxToDouble(BoxesRuntime.java:81)
    at org.apache.spark.ml.tree.CategoricalSplit.shouldGoLeft(Split.scala:107)
    at org.apache.spark.ml.tree.LearningNode.predictImpl(Node.scala:343)
    at org.apache.spark.ml.tree.impl.RandomForest$.SanonFun$findBestSplits$8(RandomForest.scala:560)
    at org.apache.spark.ml.tree.impl.RandomForest$.SanonFun$findBestSplits$8$adapted(RandomForest.scala:558)
    at org.apache.spark.ml.tree.impl.RandomForest$$Lambda$4689/342808756.apply(Unknown Source)
    at scala.collection.immutable.HashMap$HashMap1.foreach(HashMap.scala:400)
    at scala.collection.immutable.HashMap$HashMap1.foreach(HashMap.scala:728)
    at scala.collection.immutable.HashMap$HashMap1.foreach(HashMap.scala:728)
    at org.apache.spark.ml.tree.impl.RandomForest$.binSeqOps1(RandomForest.scala:558)
    at org.apache.spark.ml.tree.impl.RandomForest$.SanonFun$findBestSplits$24(RandomForest.scala:655)
    at org.apache.spark.ml.tree.impl.RandomForest$$Lambda$4688/638074108.apply(Unknown Source)
    at scala.collection.Iterator.foreach(Iterator.scala:943)
    at scala.collection.Iterator.foreach$(Iterator.scala:943)
    at org.apache.spark.InterruptibleIterator.foreach(InterruptibleIterator.scala:28)
    at org.apache.spark.ml.tree.impl.RandomForest$.SanonFun$findBestSplits$21(RandomForest.scala:655)
    at org.apache.spark.ml.tree.impl.RandomForest$$Lambda$4653/785603123.apply(Unknown Source)
    at org.apache.spark.rdd.RDD$.SanonFun$SnapPartitions$2(RDD.scala:855)
    at org.apache.spark.rdd.RDD$.SanonFun$SnapPartitions$2$adapted(RDD.scala:855)
    at org.apache.spark.rdd.RDD$$Lambda$2856/499320742.apply(Unknown Source)
    at org.apache.spark.rdd.MapPartitionsRDD.compute(MapPartitionsRDD.scala:52)
    at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:364)
    at org.apache.spark.rdd.RDD.iterator(RDD.scala:328)

```

For our experiments, we ran the program changing VM settings – we tested with 2GB, 4GB, and 6GB RAM, and 2 cores, 4 cores, and 6 cores. Tested 10-200 trees for all settings, and sometimes up to 400 for 2 cores or 4 cores. See presentation for full results and description of tests.